

Date: 29/07/2024

## WEEK 3: Hyperparameter tuning, Batch Normalization,

### Programming frameworks

Hyperparameters:

$\alpha$  → most important

{ $\beta$ } → 0.9 most default

$\beta_1, \beta_2, \epsilon 10^{-8}$

{# layers}, {# hidden units}  
{learning rate decay},  
{mini-batch size}

Priority:

1 → ①

2 → ②

3 → ③

Try random values: Don't use  
a grid (fixed values)

Date:

Fig 1.

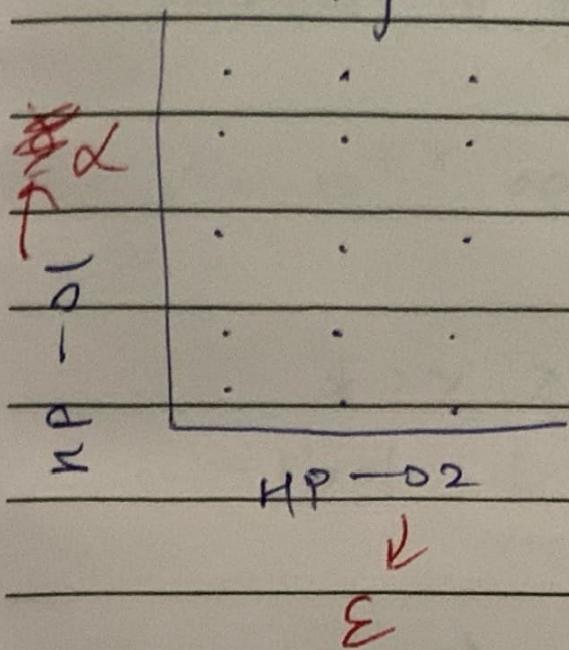
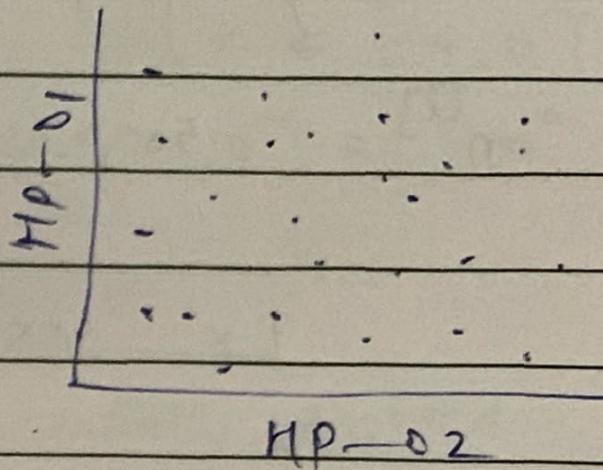


Fig 2.



In Fig 1., the values of alpha can be tried is only 5, so  $\alpha$  and  $\epsilon$  get equal importance even when importance of  $\alpha \gg \epsilon$ .

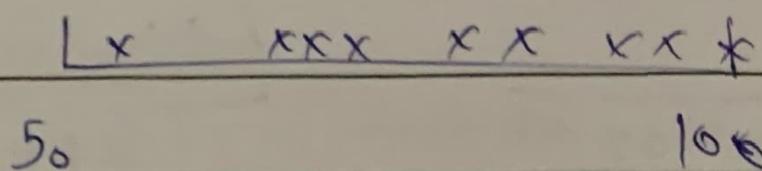
$\Rightarrow$  Coarse to fine:

We zoom in into a region where the points tend to work well. In this chosen region we find more such points similarly for best accuracy and pair-~~α~~.  $\therefore$  We can sample more densely now.

Date:

→ Picking hyperparameters at random:

$$n^{(L)} = 50, \dots, 100.$$

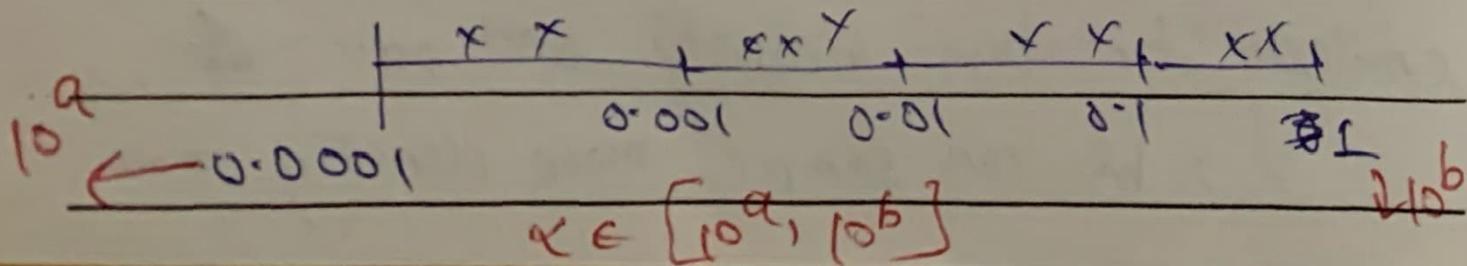
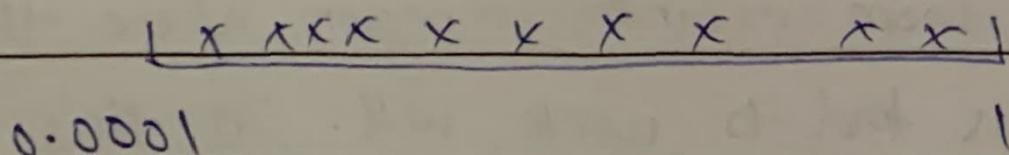


# layers  $L = 2 - 4$

Uniformly random is fine for these no. of layers.

→ Apt. scale for hyperparameters:

$$\alpha = 0.0001, \dots, 1$$



Date:

$$r = -4 * \text{np.random.rand}()$$

$$\# r \in [-4, 0]$$

$$\alpha = 10^r$$

$$\# \alpha \in 10^{-4} \dots 10^0$$

→ HiPs for exponentially weighted averages

$$\beta = 0.9 \dots 0.999$$

$\downarrow$

avg. over last 10 units

$\downarrow$

avg. over last 1000 units

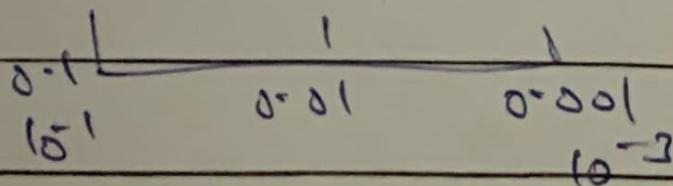
$$\lfloor x \ x x \ x \ x x x \rfloor$$

0.9

0.999

$\Gamma^B$

= 0.1 \dots 0.001



Date: \_\_\_\_\_

$$z \in [-3, -1]$$

small change in

$$1 - \beta = 10^r$$

beta close to 1

$$\beta = 1 - 10^r$$

is very sensitive,

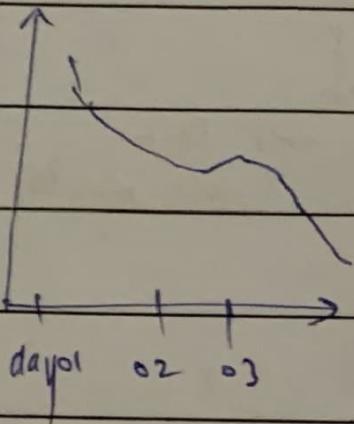
this helps in

sample more

densely when  $\beta \approx 1$

Baby sitting one  
model

Training many  
models in parallel



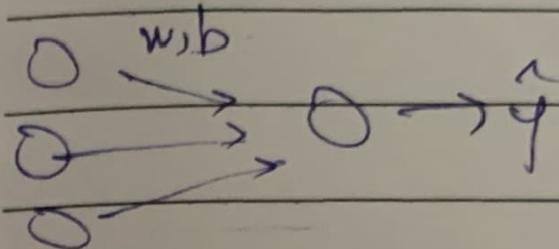
Panda

Caviar

If CPU's GPU's are not  
limited use Caviar approach.

Date:

⇒ Normalising inputs to speed up learning

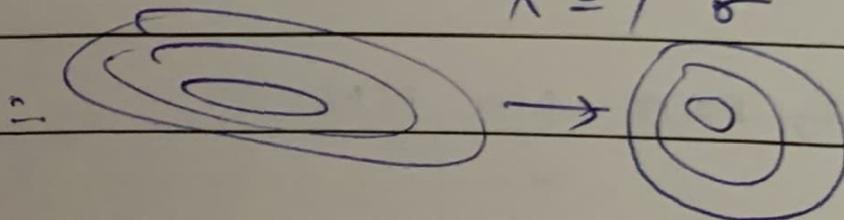


$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$x = x - \mu$$

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)})^2 \text{ element-wise}$$

$$x = / \sigma$$



⇒ Implementing a Batch norm :

Given some intermediate values

in NN

$$z^{(1)}, \dots, z^{(m)} \rightarrow z^{[e](i)}$$

$$\mu = \frac{1}{m} \sum_{i=1}^m z^{(i)}$$

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (z^{(i)} - \mu)^2$$

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

Here mean becomes  $\rightarrow 0$   
variance  $\rightarrow 1$  for  
but we don't want this for all layers

$$\therefore \Sigma^{(i)} = \gamma Z_{\text{norm}}^{(i)} + \beta$$

learnable parameters of model

this helps in setting variance and means to our choice.

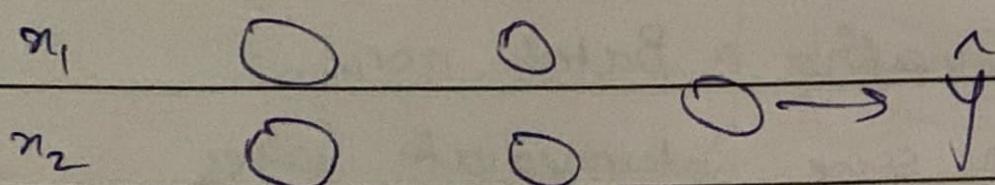
$$\text{if } \gamma = \sqrt{\sigma^2 + \epsilon}$$

and  $\beta = \mu$

then

$$\tilde{z}^{(i)} = z^{(i)}$$

→ Adding Batch norm to a NN:



The diagram illustrates a neural network layer. It starts with an input  $x$ , followed by a transformation involving weights  $w$  and bias  $b$  to produce an intermediate representation  $z$ . This is followed by a transformation involving parameters  $\beta$  and  $\gamma$  to produce the final output  $a^{[l]}$ . The final output  $a^{[l]}$  is also influenced by a normalization factor  $\sigma^2$ . Below the diagram, the text "Batch Norm (BN)" is written.

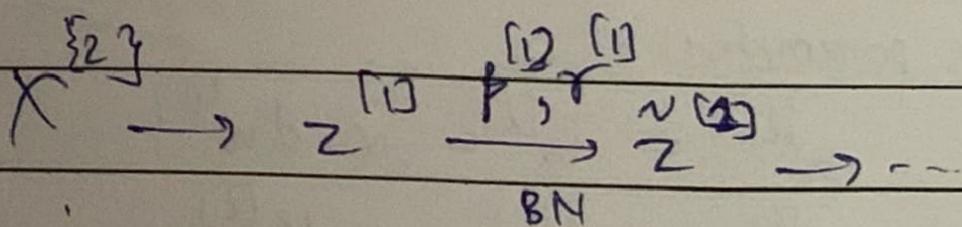
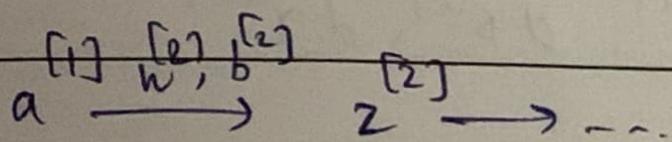
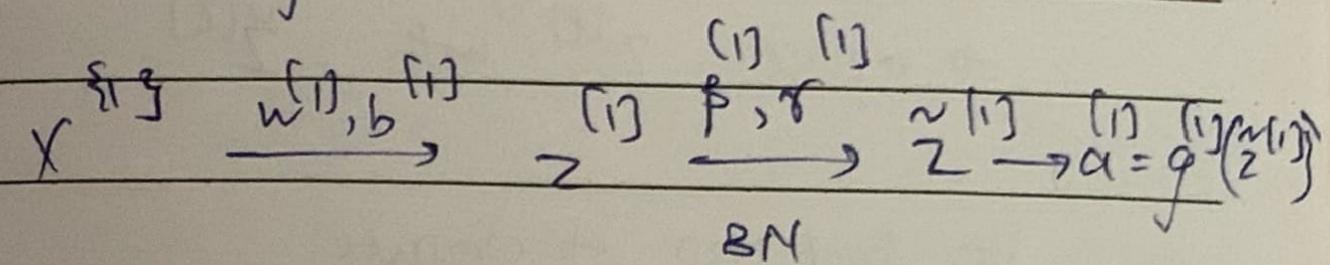
$$a^{[1]} \xrightarrow{w, b^{[2]}} z^{[2]} \xrightarrow{\beta, \gamma^{[2]}} \tilde{z}^{[2]} \xrightarrow{\alpha} a^{[1]} \rightarrow \dots$$

BN

Date:

Parameters :  $w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}$   
 $\beta^{[1]}, \gamma^{[1]}, \dots, \beta^{[L]}, \gamma^{[L]}$

→ working with mini-batches



Parameter:  $w^{[l]}, b^{[l]}, \beta^{[l]}, \gamma^{[l]} \xrightarrow{\text{oder}} (n^{[l]}, \cdot)$

During BN,  $b^{(x)}$  can be

eliminated or be set permanently to zero.  
because batch norm is equal to the mean

because batch-norm zeroes out the mean

→ No point in adding everywhere: if  
is going to be multiplied

Date:

Implementing grad. descent.

for  $t = 1 \dots \text{num Minibatches}$

Compute forward prop on  $X^{\{t\}}$

In each hidden layer, use BN

to replace  $Z^{[l]}$  with  $\hat{Z}^{[l]}$

Use back prop to compute

$dw^{[l]}$ ,  $d\beta^{[l]}$ ,  $d\gamma^{[l]}$

Update parameters

$$w^{[l]} := w^{[l]} - \alpha dw^{[l]}$$

$$\beta^{[l]} := \beta^{[l]} - \alpha d\beta^{[l]}$$

$$\gamma^{[l]} := \gamma^{[l]} - \alpha d\gamma^{[l]}$$

works with momentum, RMS prop, Adam.

Date:

→ When the distribution of the input set changes, the parameters learnt further change, and ∴ the deep NN suffers a covariate shift, causing less accuracy and the cost function might not minimise well.

∴ If batch norm is applied, the distributions are normalised helping the model to learn the parameters in a better way.

→ Batch norm as regularization :

① Each mini-batch is scaled by mean/variance computed on just that mini-batch

② This adds some noise to the values

$z^{[l]}$  within that mini-batch. So similar

to dropouts it adds some noise to each hidden layer's activation.

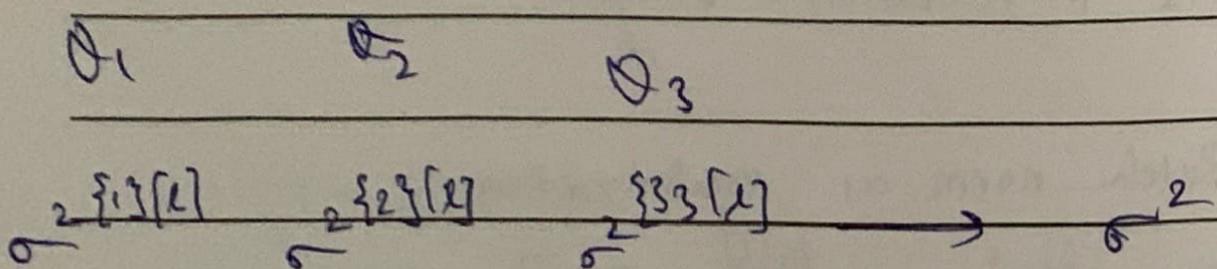
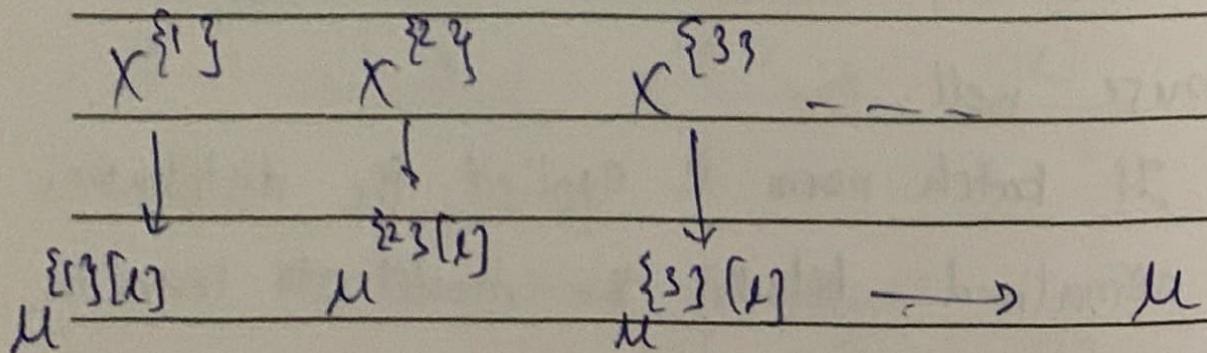
③ This has a slight regularization effect  
batch-norm + dropout

= powerful regularization

Date:

→ Batch-norm at test time:  
(Running average)

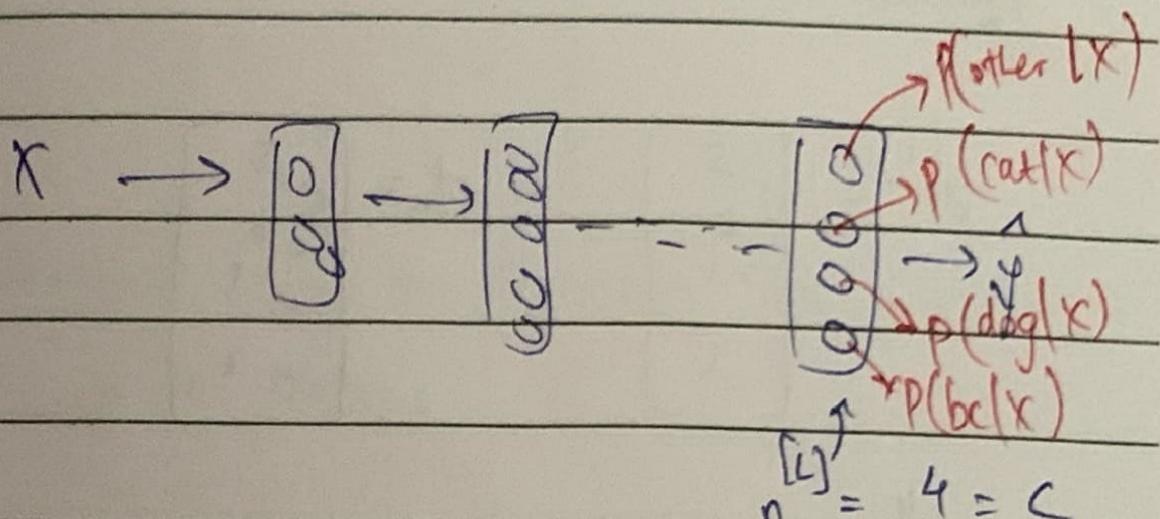
$\mu, \sigma^2$ : estimate using exponentially  
weighted average (across mini-batches)



Date:

→ Softman regression: (generalization of LR  
with ~~two~~ or more  
classes)

$$C = \# \text{ classes} = 4 \quad (0, \dots, 3)$$



$$\mathbf{y} \in \{0, 1\}$$

Softman layer:

$$z^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]} \quad (4, 1)$$

Activation function:

$$t = e^{z^{[L]}} \quad (4, 1)$$

Date:

$$a^{[L]} = \frac{z^{[L]}}{\sum_{i=1}^4 t_i}$$

$$a_i^{[L]} = \frac{t_i}{\sum_{i=1}^4 t_i}$$

Example:

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}$$

$$\sum_{i=1}^4 t_i = 176.3$$

$$a^{[L]} = \underline{\langle t \rangle} \rightarrow \text{as a vector}$$

$$176.3$$

individually;

$$\text{for } i=1 \quad e^5 / 176.3$$

$$= 2 \quad e^2 / 176.3$$

$$= 3 \quad e^{-1} / 176.3$$

$$= 4 \quad e^3 / 176.3$$

{ These  
sum to  
1.

Soft TMAX:

Date: 30/07/2024

Softmax example: Shown in video.

"hard man"

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \xrightarrow{\text{gives the highest value.}}$$

Softmax regression generalises logistic regression to C classes.

If  $C=2$ , softmax = logistic regression

Loss function:

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

$$a^{[C]} = \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$C=4$  here

Date:

$$\mathcal{L}(\hat{y}, y) = - \sum_{i=1}^C y_i \log(\hat{y}_i)$$

$$= -y_2 \log \hat{y}_2 = -\log \hat{y}_2$$

To make  $-\log \hat{y}_2 \downarrow$

make  $\hat{y}_2 \uparrow$

$$J(w^{(1)}, b^{(1)}, \dots) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(y^{(i)}, \hat{y}^{(i)})$$

$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$$

$$= (C, m)$$

$$\hat{Y} = [\hat{y}^{(1)} \ \hat{y}^{(2)} \ \dots \ \hat{y}^{(m)}]$$

Date:

Softmax layer outputs

$$z^{[L]} \rightarrow a^{[L]} = \hat{y} \rightarrow L(\hat{y}, y)$$

Back prop:

$$dz^{[L]} = \hat{y} - y$$